

# Человек-вектор и подъем конверсии: как мы сделали шаг к Federated Learning

Анастасия Семенова  
Иван Снегирев  
Артём Просветов



**HighLoad++**  
Весна 2021



# О нас



Анастасия Семенова

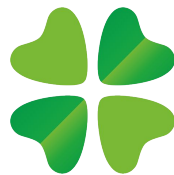
- Ведущий Data Scientist (CleverDATA)
- Веду семинары в ВШЭ



Артём Просветов

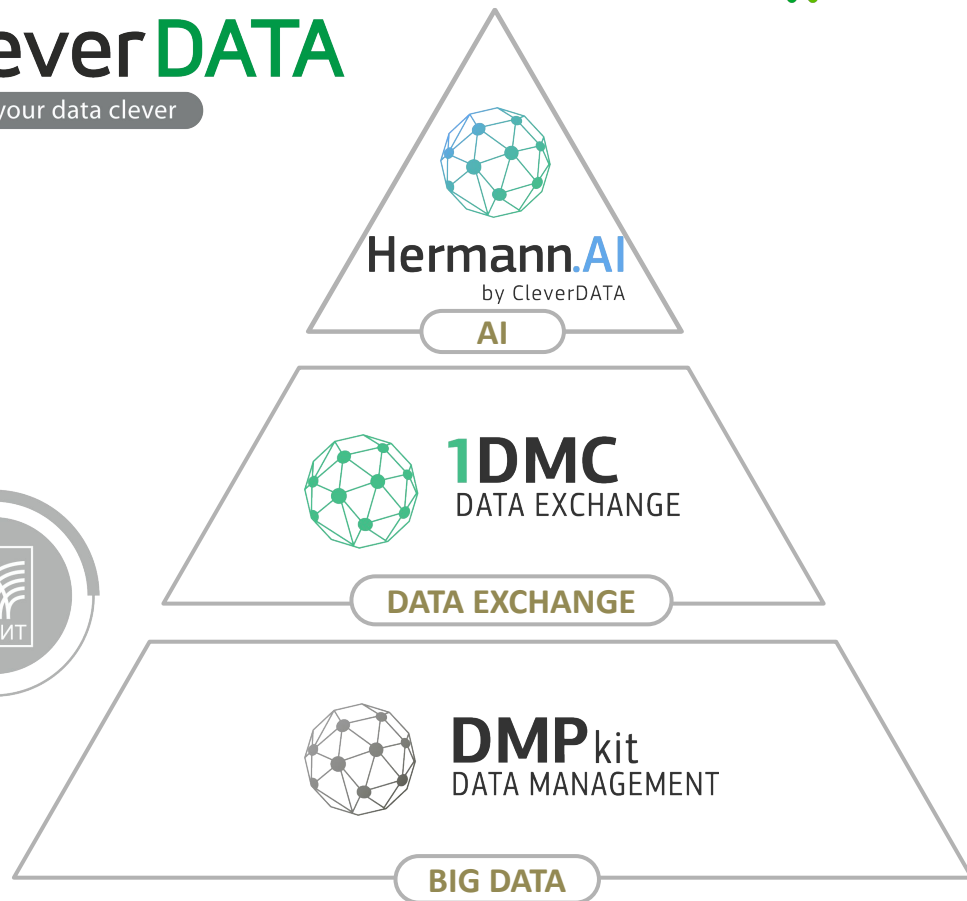
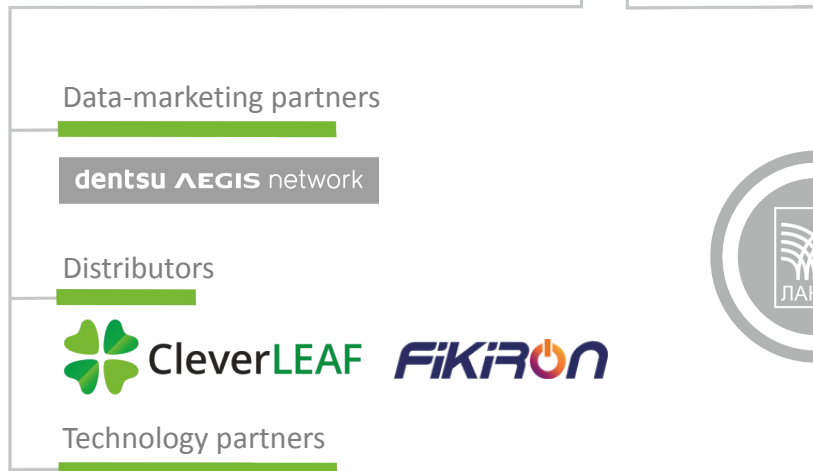
- Chief Data Scientist (LANIT/CleverDATA)
- Ведущий математик ИКИ РАН
- Веду лекции в ВШЭ
- к. ф.-м. н.

CLEVERDATA –  
ТЕХНОЛОГИЧЕСКИЙ  
ПРОВАЙДЕР РЕШЕНИЙ ДЛЯ  
МОНЕТИЗАЦИИ ДАННЫХ

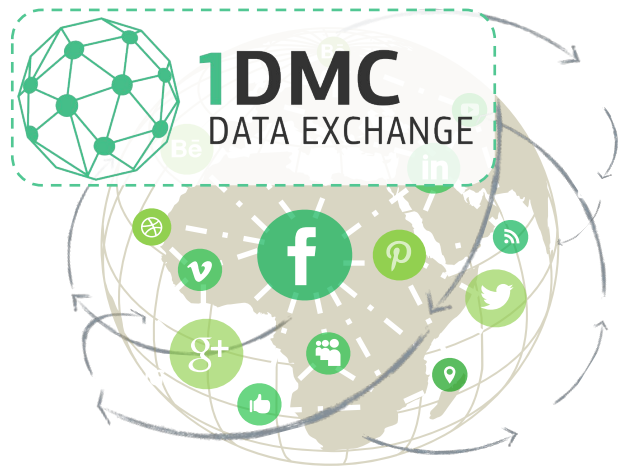


# CleverDATA

make your data clever



# ЕДИНАЯ ТОЧКА ИНТЕГРАЦИИ ПОСТАВЩИКОВ И ПОТРЕБИТЕЛЕЙ ВНЕШНИХ ДАННЫХ



**1DMC БИРЖА ДАННЫХ** – площадка, объединяющая поставщиков и потребителей данных. Доступный объем данных сопоставим с крупнейшими в России сайтами.

ПОСТАВЩИКИ ДАННЫХ	20+
ИСТОЧНИКИ ДАННЫХ	9000+
АТРИБУТЫ	30000+
СУТОЧНАЯ АУДИТОРИЯ	85M
РЕКЛАМНЫЕ СИСТЕМЫ &	10+

## ПРИМЕНЕНИЕ ВНЕШНИХ ДАННЫХ

Таргетированная онлайн реклама	Медийные охватные кампании	
Динамический контент	Look-alike	Повторные продажи
Скоринг	Антифрод	Brand awareness и точность таргетинга

# Оглавление

- CleverDATA
- Задачи
- 2 проблемы
  - Текущее представление профиля человека и вопросы к сегментации
  - Что такое FL и почему он нам интересен
- Решение – эмбединги!
- Подход 1: Текст vs Граф
- Подход 2: Time Encoder

# Задачи

1. Lookalike (LaL)
2. Вероятность целевого действия

# Проблемы

**человек**  $\mapsto$  (0, 0, 0, ... 0, 1, 0, 0, ... 0, 1, 1, 0, 0, ... 0, 0)

интерес к велосипедам  $\swarrow$

читает про аквариумных рыбок  $\swarrow$

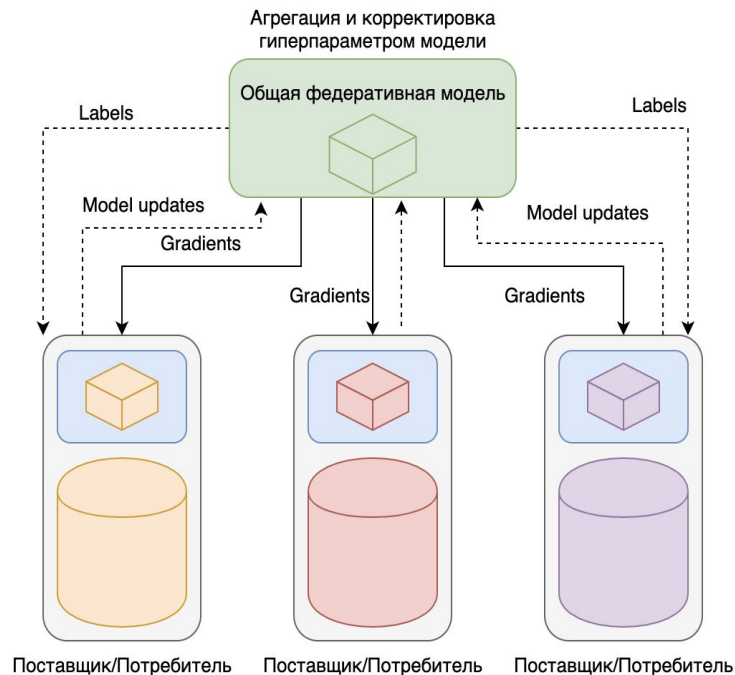
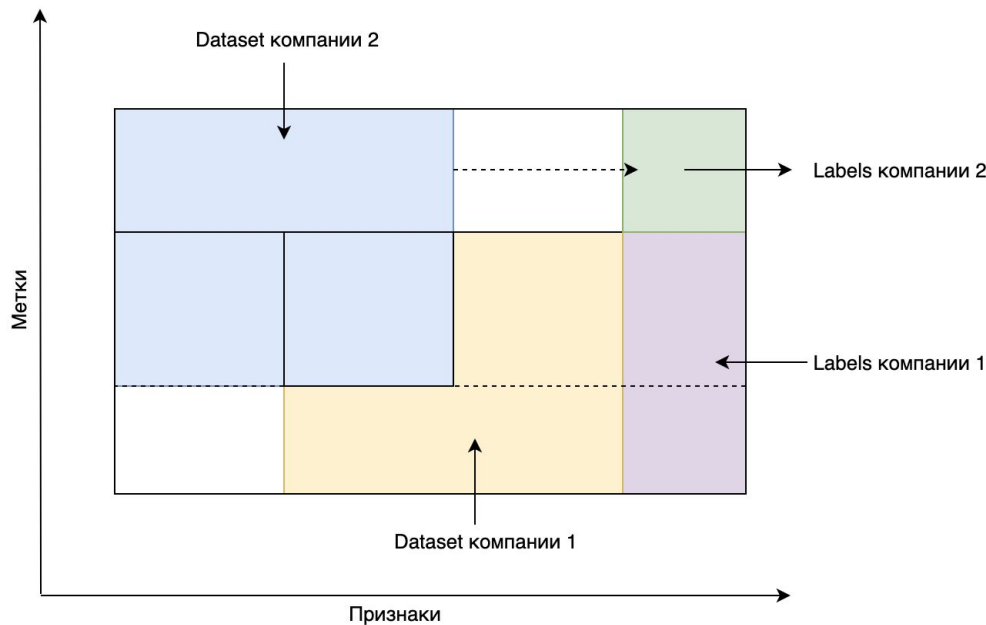
смотрит теннис  $\nearrow$

# Проблемы

1. Слишком **длинный вектор** профиля
2. Данные разрежены
3. Высокая размерность пространства представлений
4. Хотим **Federated Learning**



# Federated Learning



# 1. Lookalike

**Lookalike (LaL)** – это таргетинг, при котором рекламные материалы показываются тем пользователям, которые по поведенческим характеристикам похожи на текущую аудиторию ресурса.

**Clickstream** – сегмент целевых пользователей (целевой класс)

**Repr\_sample** – сегмент случайных пользователей

- **tf-idf**
- **LogisticRegression** / boosting / etc.

## 2. Предсказание осуществления целевого действия

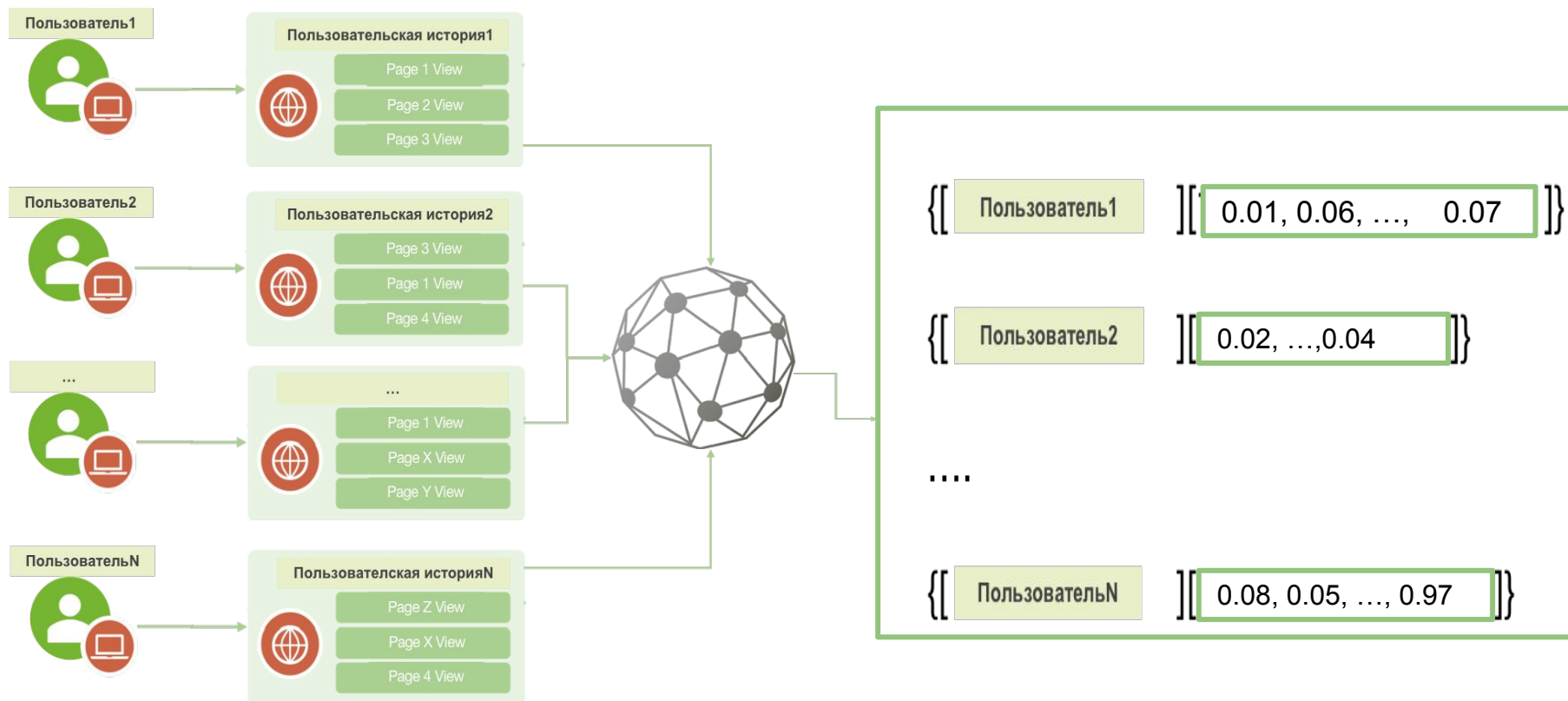
Выбирается целевое действие на сайте клиента DMP.

Задача **бинарной классификации**: те, кто осуществил целевое действие (target) и кто – нет (not\_target).

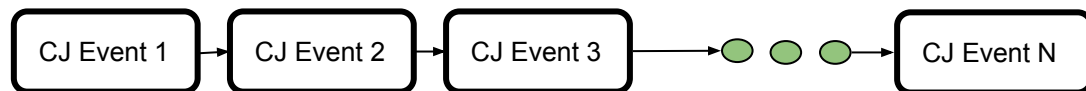
# Проблемы

1. Слишком длинный вектор профиля
2. Данные разрежены
3. Высокая размерность пространства представлений
4. Хотим Federated Learning

**Решение: Embeddings!**

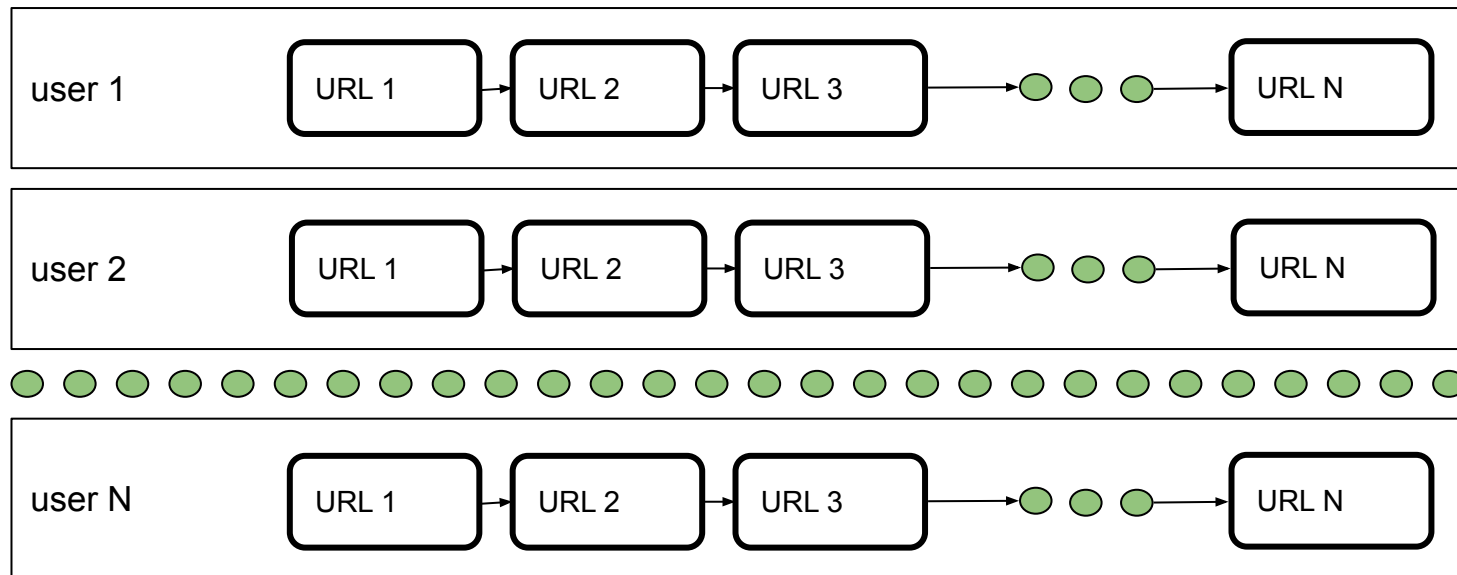


# Customer Journey (далее CJ)



- Ссылка, с которой перешел пользователь
- Время события
- Уникальные идентификаторы пользователя
- UserAgent
- Cookies
- Headers
- Params

# Подготовка данных



# История первая

Текстовая векторизация



Мир меняется... А ссылки?

человек-  
текст

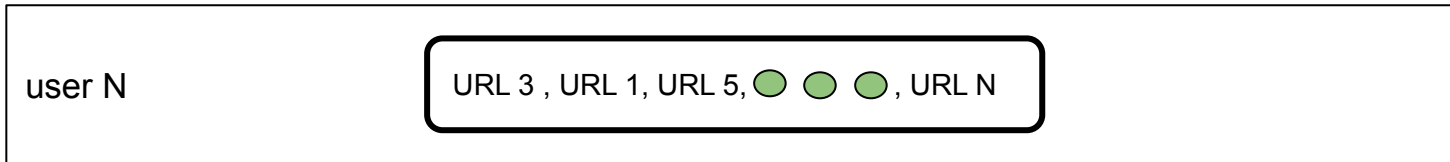
Мир меняется... А ссылки?

Был выбран **FastText** с буквенными **n-граммами**.

Мир меняется... А ссылки?

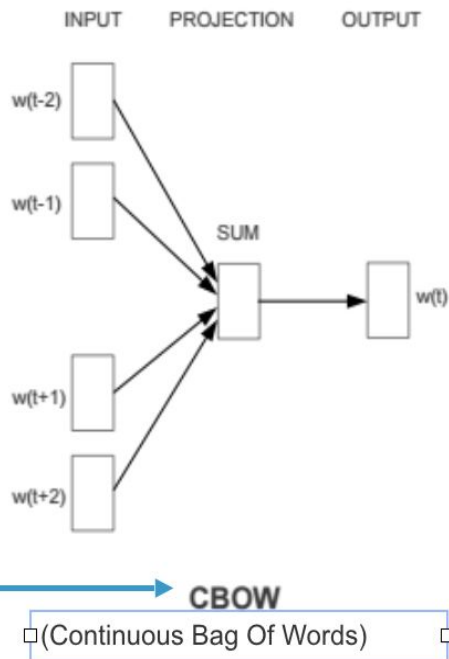
Был выбран FastText с буквенными n-граммами.

**Пользователи** = последовательности – это **документ**  
Ссылка – это **токен**

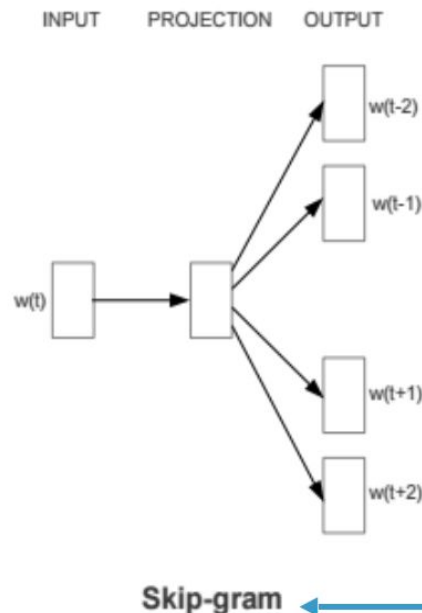


ЧЕЛОВЕК-  
ТЕКСТ

Given a set of  
(neighboring) words,  
**guess single words**  
that potentially occur  
along with this set of  
words.



or



**Guess potential  
neighboring  
words** based on  
the single word  
being analyzed.

<https://derekchia.com/an-implementation-guide-to-word2vec-using-numpy-and-google-sheets/>

<https://mailsgun.ru/deep-learning-vs-common-sense-разрабатываем-чат-бота/>

человек-  
текст

FastText от Gensim

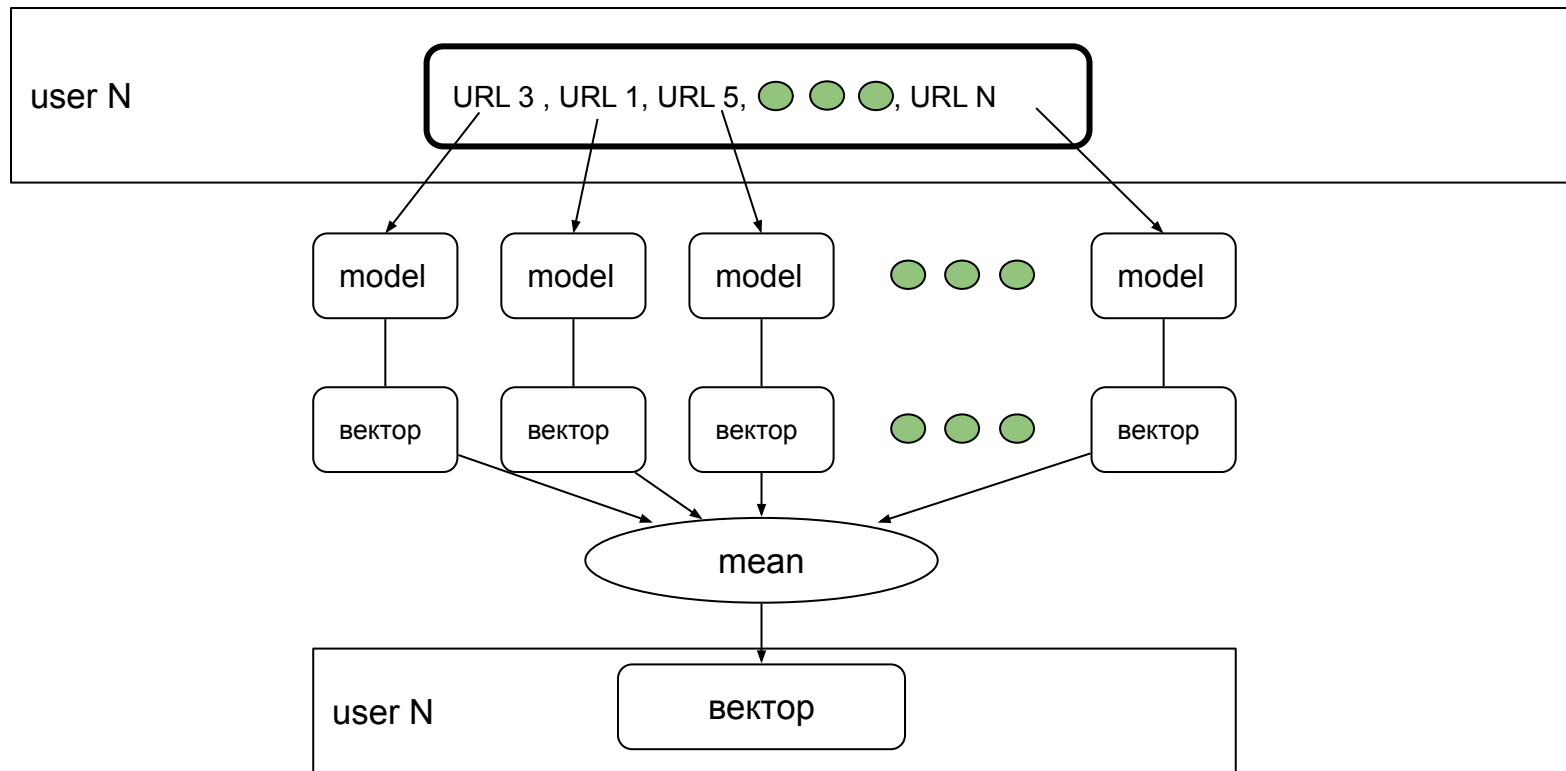
**Skip-Gram**

**dim = 128**

**min count token = 5**

**окно = 5**

ЧЕЛОВЕК-  
ТЕКСТ



человек-  
текст

**blabla**\_url.com *нет* в словаре

человек-  
текст

**blabla\_url.com** *нет* в словаре

**bla**\_url.com *есть* в словаре



человек-  
текст

**blabla\_url.com** *нет* в словаре

**bla\_url.com** *есть* в словаре

n-gramm => :- ) => **есть** вектор для  
**blabla\_url.com!**

# Плюсы

- Реализация Gensim позволяет обучаться на генераторах
- Метрики качества выше, чем у классического подхода
- Позволяет работать с незнакомыми токенами
- Большой прирост полноты относительно классического подхода

человек-  
текст

Подход	ROC AUC	F1-мера для целевого класса	Precision для целевого класса	Recall для целевого класса
classic	0.633	0.053	0.056	0.051
text_embs	0.664 (+4.9 %)	0.056 (+5.7 %)	0.032 (-42 %)	0.224 (+337 %)

человек-  
текст

Всё здорово, но чего-то не хватает

??? человек - не текст, человек - граф!

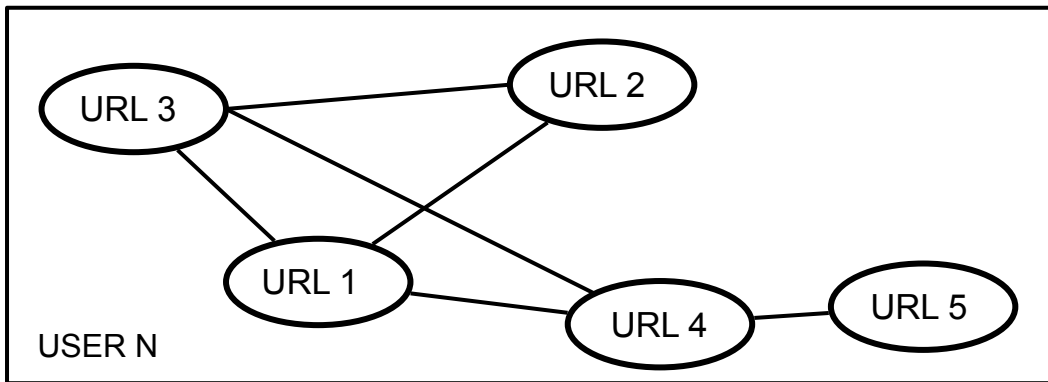
# История вторая

Графовая векторизация

# Структура переходов пользователей интернета важна!

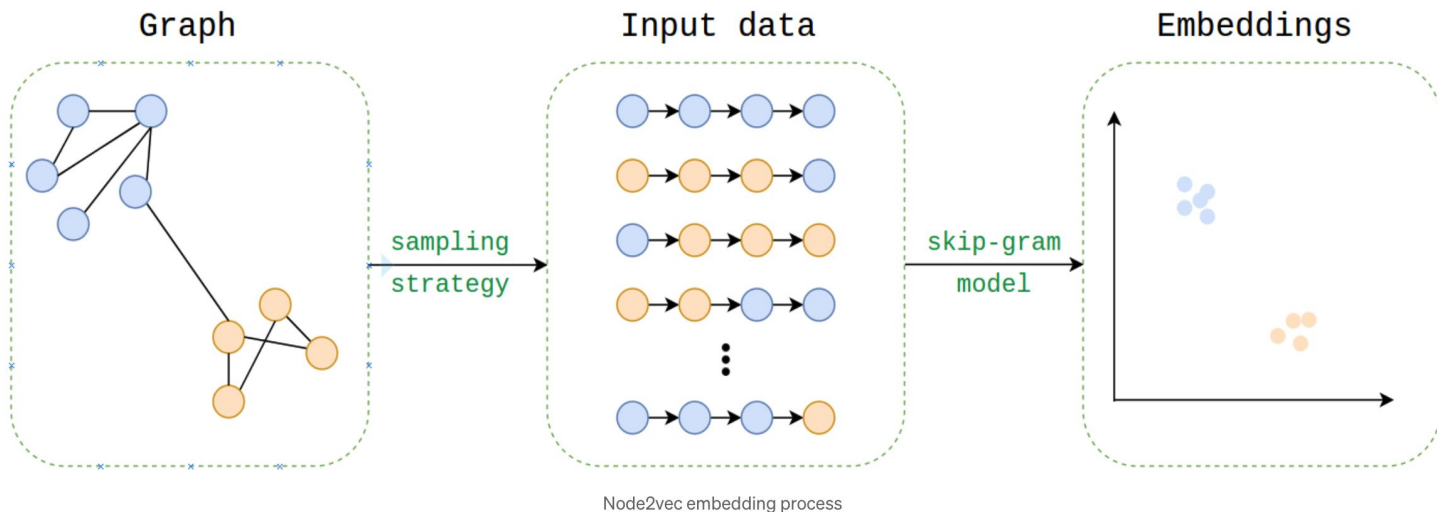
Нижняя оценка размера графа на месячных данных составляет 160 000 доменных узлов.

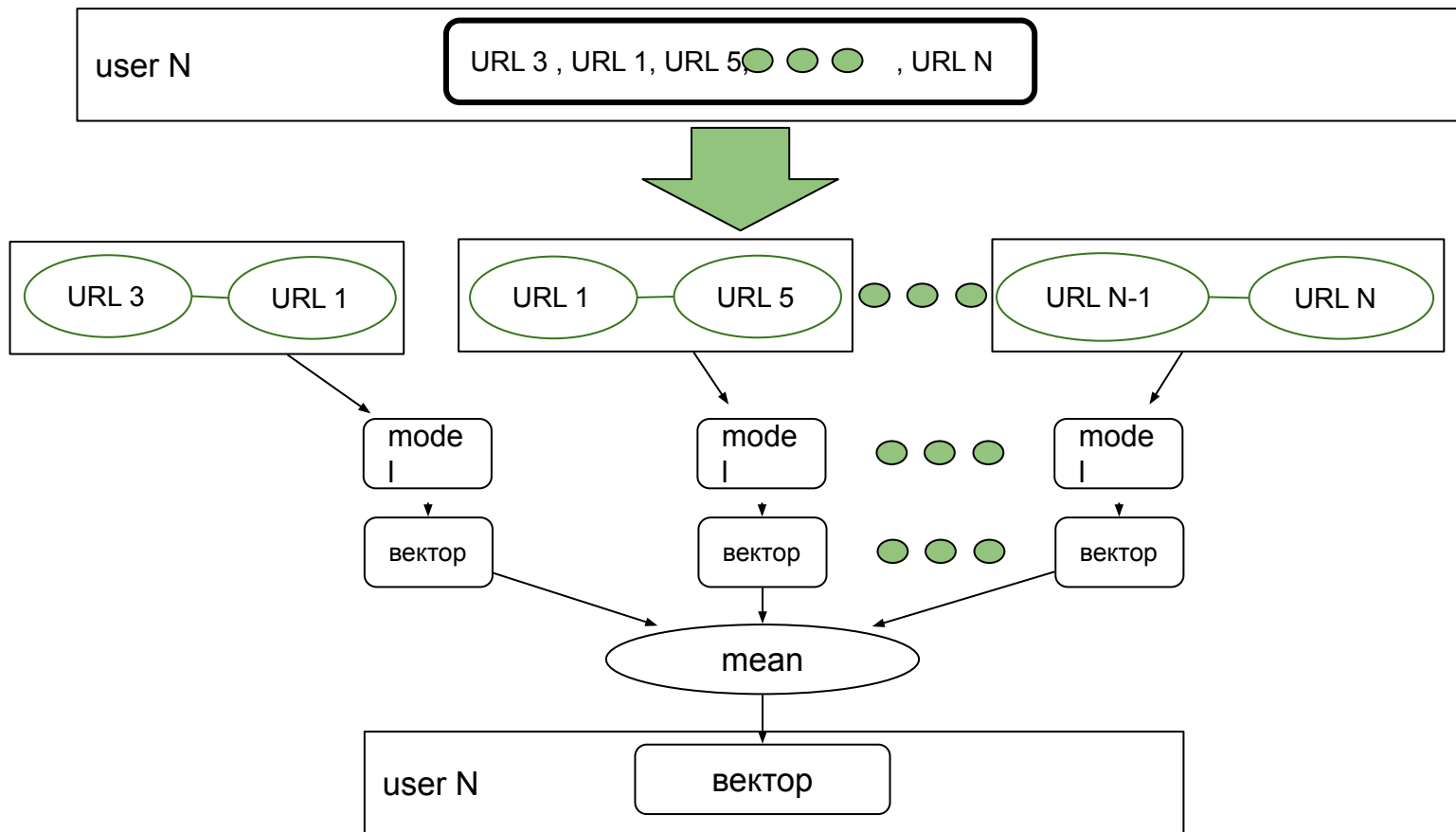
При этом мы постоянно подключаем новых поставщиков данных.



# Node2Vec

отображает соседние узлы в графе так, чтобы и в признаковом пространстве они были близки







# Плюсы

- Хорошая **точность** для уже известных токенов
- Время обучения не зависит от количества данных, только от размера графа
- Метрика **полноты выше**, чем у классического подхода

# Минусы

- Нельзя рассчитать вектор для **неизвестного токена**
- Если у пользователя нет ни одного известного токена, то приходится брать среднее среди всех пользователей
- Рассчитываем вектор только для ребра, а не всей последовательности
- В целом показатели **качества ниже**, чем у классического подхода

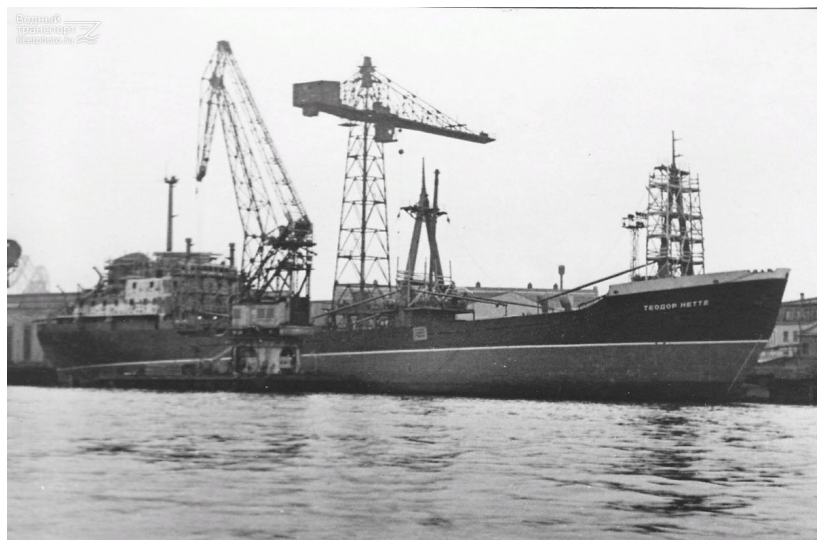
Подход	ROC AUC	F1-мера для целевого класса	Precision для целевого класса	Recall для целевого класса
classic	0.633	0.053	0.056	0.051
graph_embs dim 128	0.627 (-0.9 %)	0.039 (-27.4 %)	0.023 (-58.7 %)	0.122 (+137.5 %)

?????? Метрики Precision и Recall получены в результате перебора порога бинаризации для максимизации F1-меры для целевого класса

Человек — текст

или

Человек — граф?



# Теорема Кондорсе!

# Теорема Кондорсе!

или по-простому: пробуем **объединить** подходы

**текст (dim 128) + граф (dim 128) = emb (dim 256)**

Подход	ROC AUC	F1 мера для целевого класса	Precision для целевого класса	Recall для целевого класса
classic	0.632 (0 %)	0.053 (0 %)	0.056 (0 %)	0.051 (0 %)
text_embs dim 128	0.664 (+4.9 %)	0.056 (+5.7 %)	0.032 (-42 %)	0.224 (+337 %)
graph_embs dim 128	0.627 (-0.9 %)	0.039 (-27.4 %)	0.023 (-58.7 %)	0.122 (+137.5 %)
embs_union dim 256	<b>0.679 (+7.3 %)</b>	<b>0.062 (+16.6 %)</b>	0.043 (-23 %)	0.115 (+125 %)

Метрики Precision и Recall получены в результате перебора порога бинаризации для максимизации F1 мера для целевого класса



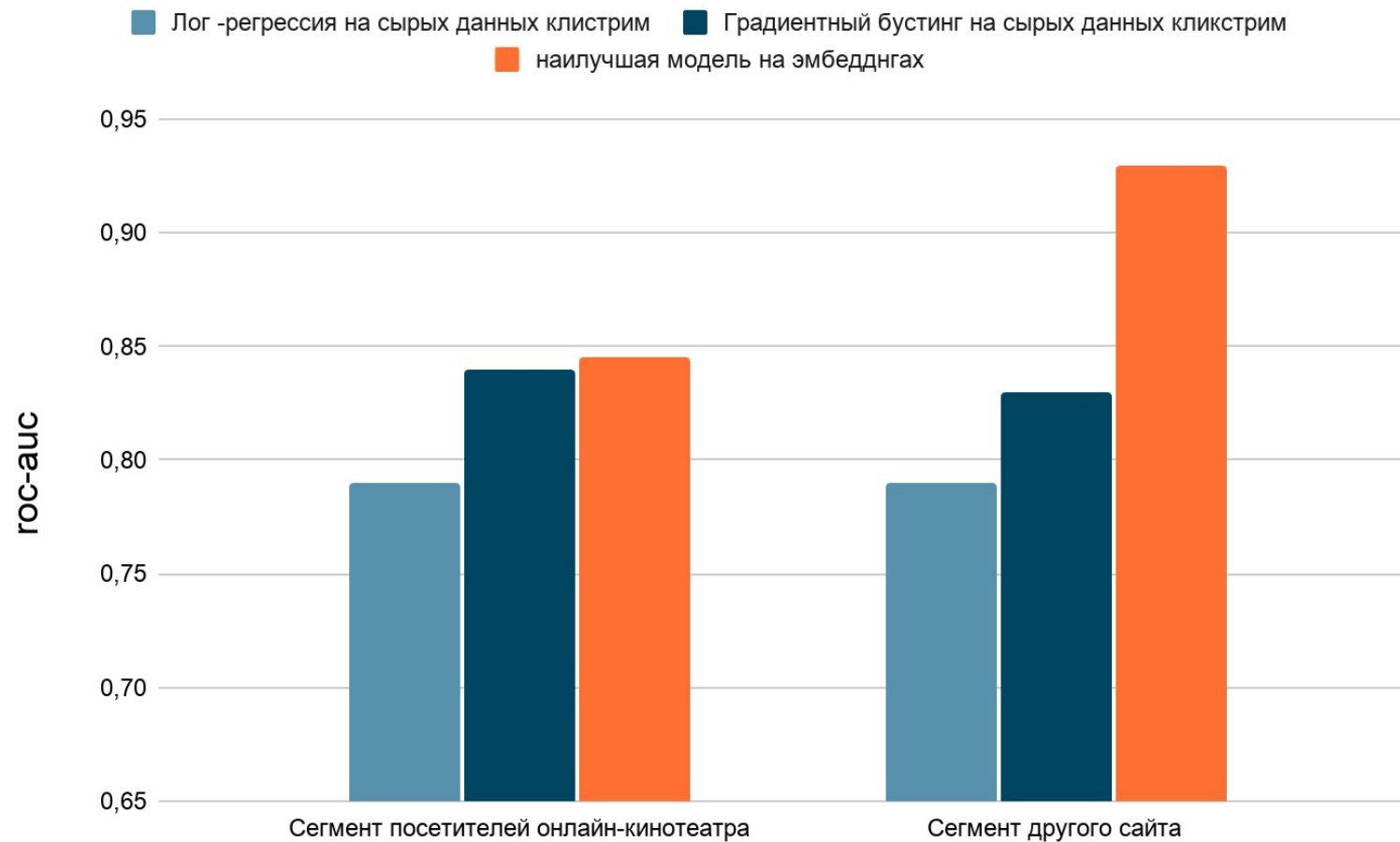
# Кейс

Онлайн-кинотеатр

Неглубокая **полносвязная сеть**

В **20 раз лучше**, чем случайный  
предиктор





Что можно **ИСПОЛЬЗОВАТЬ**  
еще?

# Что можно использовать еще?

- Ссылка, с которой перешел пользователь
- Время события
- Уникальные идентификаторы пользователя
- UserAgent
- Cookies
- Headers
- Params
-

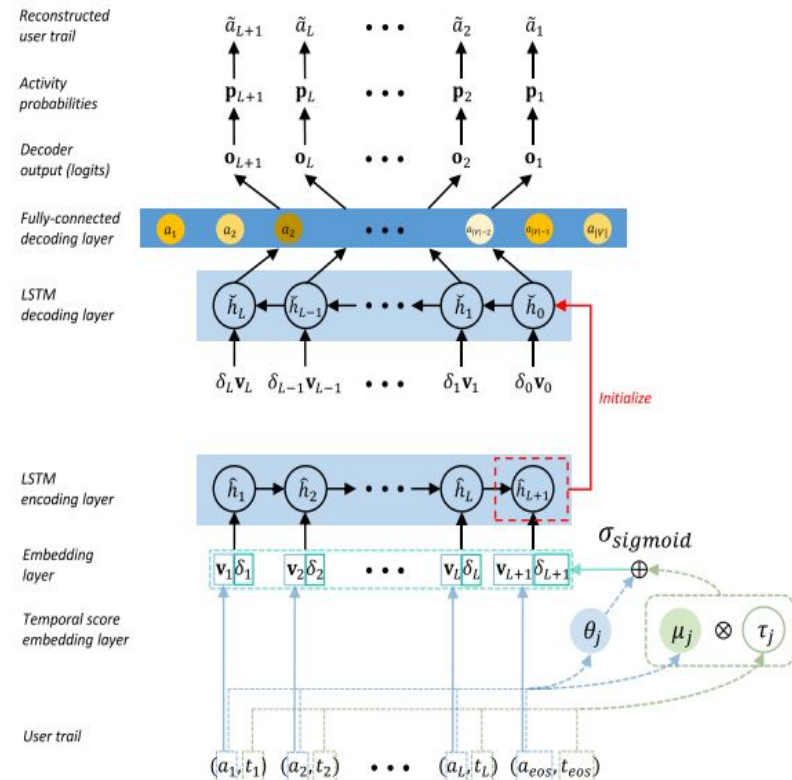
# Что можно использовать еще?

- Ссылка, с которой перешел пользователь
- **Время события**
  - Уникальные идентификаторы пользователя
- UserAgent
- Cookies
- Headers
- Params
-

# История третья

TimeEncoder

**Sequence-to-sequence модель** – это модель, принимающая на вход последовательность и возвращающая другую (такую же) последовательность элементов.

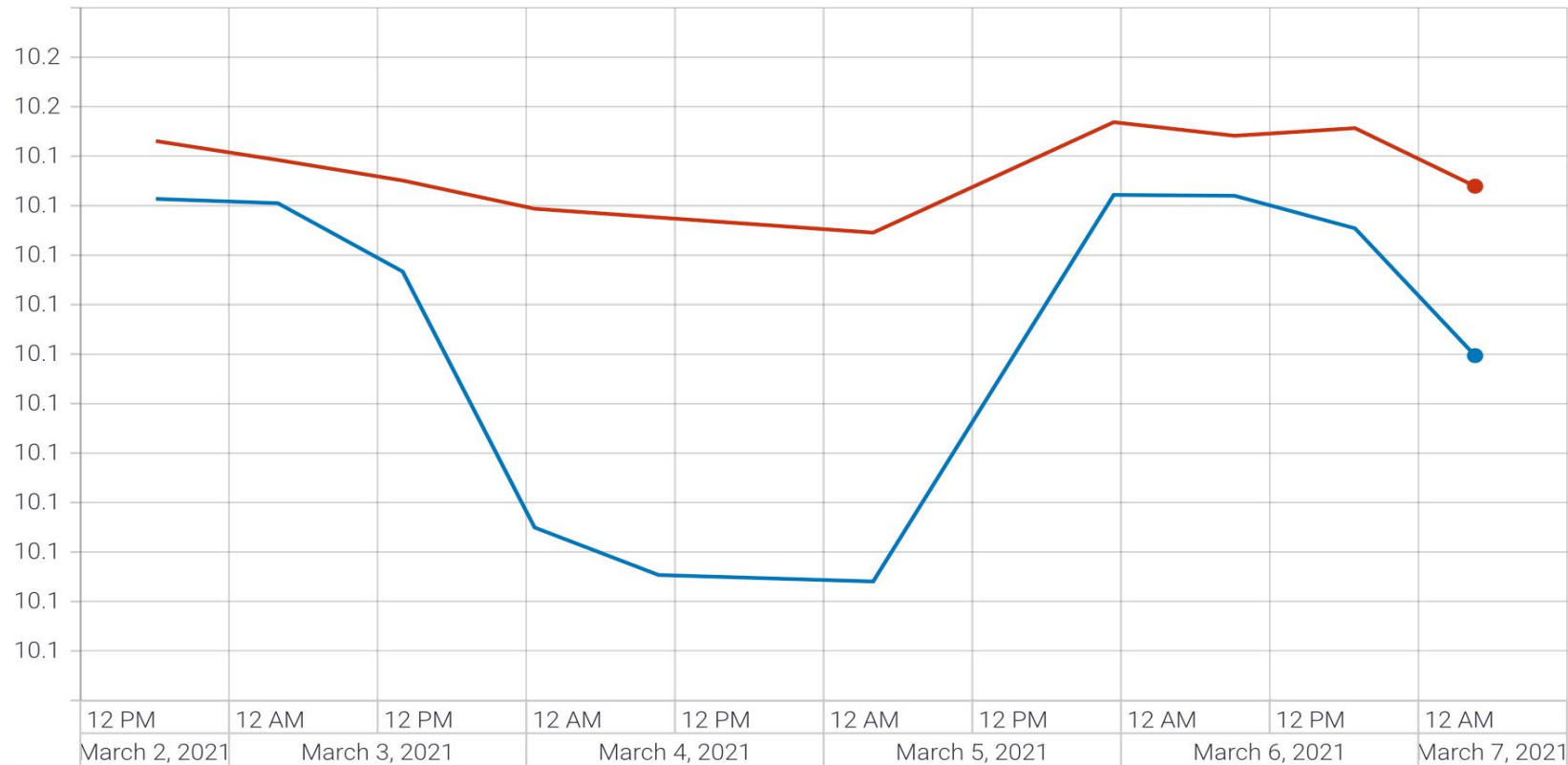


# Параметры модели

- Размерность LSTMs, Embedding = 100
- Максимальная длина последовательности 64 (95% сессий имеют 61 и меньше событий в сессии)
- Количество эпох 10
- Оптимайзер Adam ( $lr = 0.001$ )
- Изменение CrossEntropyLoss в процессе обучения



epoch  
tag: loss/epoch



# Векторизатор

- Для процесса получения векторов необходимо **удалить** декодирующую часть (**DECODER**) и оставить только кодирующую (ENCODER)
- **Усредненный** вектор выходов **LSTM**-слоя

# Как тестировали эмбеддинги?

# 1 - Lookalike

**Сравнение** с нашим подходом, объединяющим  
**Текстовые** и **Графовые эмбединги**

Среди всех model ищется соответствующий метрике  
**максимум**, кроме случаев переобучения

# 1 - Lookalike

Сравнение с нашим подходом, объединяющим  
**Текстовые** и **Графовые** эмбединги

LaL CAR\_X

Превосходство TimeEncoder над подходом graph\_text:

- roc\_auc **+22.5%**
- f1\_score clickstream **+75%**
- f1\_score repr\_sample **+1.42%**

# 1 - Lookalike

Сравнение с нашим подходом, объединяющим  
**Текстовые** и **Графовые** эмбединги

LaL CAR\_Y

Превосходство TimeEncoder над подходом graph\_text:

- roc\_auc **+8.23%**
- f1\_score clickstream **+34.7%**
- f1\_score repr\_sample **+2%**

# 1 - Lookalike

Сравнение с нашим подходом, объединяющим  
**Текстовые** и **Графовые** эмбединги

LaL APPARTS

Превосходство TimeEncoder над подходом graph\_text:

- roc\_auc **+6.5%**
- f1\_score clickstream **+25%**
- f1\_score repr\_sample **+2%**

# 1 - Lookalike

Сравнение с нашим подходом, объединяющим  
**Текстовые** и **Графовые** эмбединги

LaL AVIA

Превосходство TimeEncoder над подходом graph\_text:

- roc\_auc **+0.62%**
- f1\_score clickstream **-3%**
- f1\_score repr\_sample **-0.12%**



# 1 - Lookalike

Сравнение с нашим подходом, объединяющим  
**Текстовые** и **Графовые** эмбединги

**Превосходство** TimeEncoder над подходом graph\_text (**в среднем**):

- roc\_auc **+9.46%**
- f1\_score clickstream **+32.93%**
- f1\_score repr\_sample **+1.325%**

## 2 - Предсказание осуществления целевого действия

Превосходство TimeEncoder над подходом graph\_text:

- roc\_auc **+48%**
- f1-score target **+283%**
- f1-score not\_target **+168%**
- f1-score\_max for target **+281%**

# Дальнейшие шаги

- Изучить различные способы представления **времени**
- ...
- А работает ли это на **других** данных?

# Дальнейшие шаги

- Изучить различные способы представления времени
- ...
- А работает ли это на **других** данных?

**Clickstream**

**ОФД**

Доставка - Стандартная доставка в течение 2-3 ...	0.00	2020-09- 19 00:00:00
Бальзам для губ Lip Juicer Малина, свекла и им...	632.00	2020-09- 19 00:00:00
Скраб для тела French Grape Seed	1592.00	2020-09- 19 00:00:00
Йогурт для тела «Миндальное молочко»	792.00	2020-09- 19 00:00:00
Пилинг для тела «Клубника»	712.00	2020-09- 19 00:00:00

# Данные ОФД

Текущий пайплайн работает хорошо

<https://m.youtube.com/watch?v=R304Xy5eTso&index=105>

<https://habr.com/ru/company/lanit/blog/462959/>

# Данные ОФД

Текущий пайплайн работает хорошо

<https://m.youtube.com/watch?v=R304Xy5eTso&index=105>

<https://habr.com/ru/company/lanit/blog/462959/>

**Хотим:**

- Менее разреженные представления
- Неинтерпретируемые
- Вектор короче

# Данные ОФД

**Пилот с другим онлайн-телевидением:**

Задача холодного старта рекомендательной системы



# Данные ОФД

Пилот:

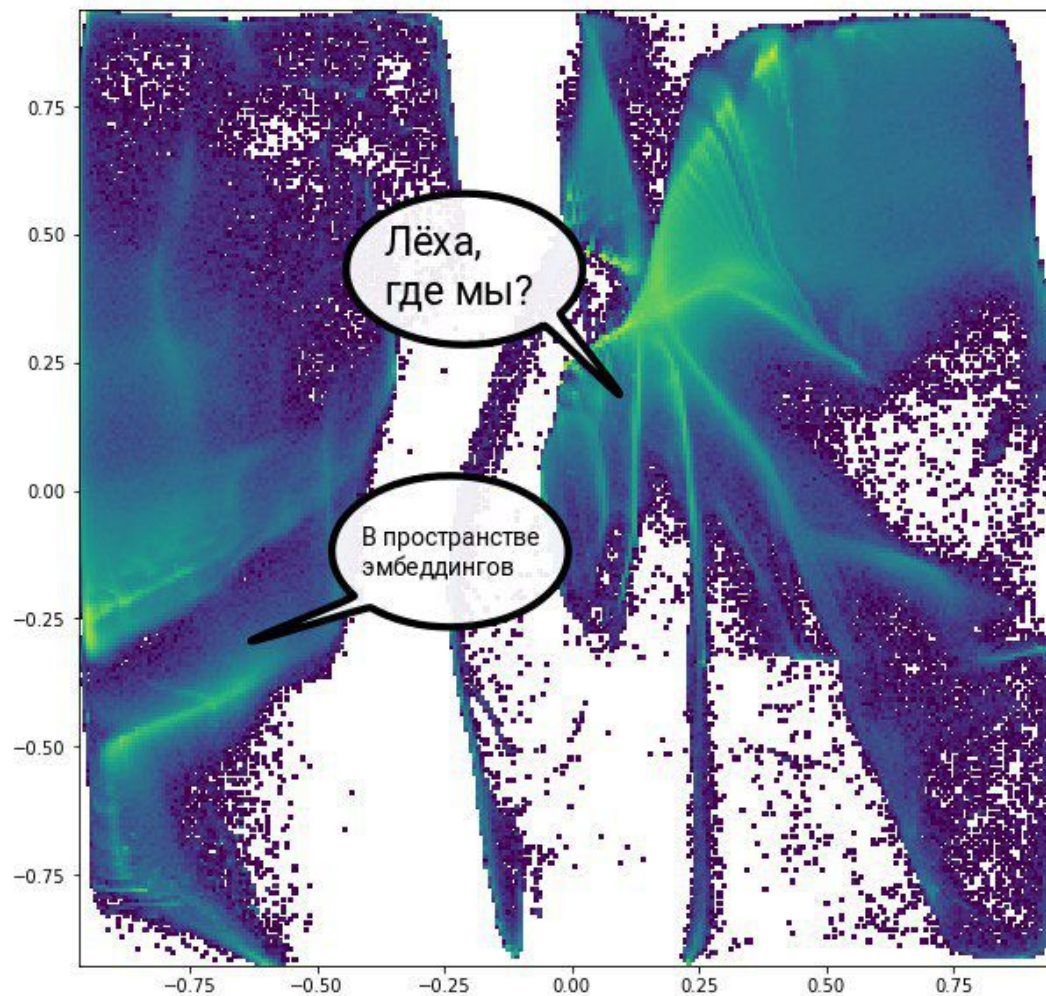
Проблема **холодного старта**  
рекомендательной системы

**Успех!** Ждём пресс-релиз



# Эмбединги на данных ОФД

	Эмбединги
ROC-AUC	- 4%
Average Precision	+ 92%
NDCG	+ 40%
NDCG@50	+ 34%



# Выводы

Правильно **подготовленные эмбединги** способны на многое!

**Подготовить наилучшим образом эмбединги помогут:**

- **статьи** других исследователей
- добавление важной информации (например, **времени**)
- подходы из **различных** направлений
- **эксперименты** на различных задачах

Что-то **новое** может быть улучшенным **старым**

А может быть **совсем новым** :-)

Будьте **открыты** идеям и экспериментам!



[https://t.me/samy\\_1010](https://t.me/samy_1010)



[https://t.me/Art\\_pr0](https://t.me/Art_pr0)

# Спасибо!

Анастасия Семенова  
Иван Снегирев  
Артём Просветов



**HighLoad++**  
Весна 2021

